

CS238 - Project 1 README

Shreya S Ramanujam

October 15, 2025

1 Description of Algorithm Used

For this project, I used the K2 algorithm followed by a greedy local search (add/remove edges greedily) to further optimize the graph we get from K2. The pseudocode for the algorithm is detailed below:

Algorithm 1: Bayesian Network Structure Learning using K2 and Greedy Hill Climbing

Input: Data matrix D with n discrete variables

Output: Directed acyclic graph G

Initialization:

Initialize an empty directed graph G with n nodes.

Set a variable order $[1, 2, \dots, n]$.

Compute initial score of empty graph $S \leftarrow \text{BayesianScore}(G, D)$.

K2 Edge Addition Phase:

for $i \leftarrow 1$ to n do

 for $j \leftarrow i + 1$ to n do

 Add edge $(X_i \rightarrow X_j)$ to G .

 Compute $S_{\text{new}} \leftarrow \text{BayesianScore}(G, D)$.

 if $S_{\text{new}} > S$ then

$S \leftarrow S_{\text{new}}$

// Keep the edge

 else

 Remove edge $(X_i \rightarrow X_j)$

// Revert if score decreases

Greedy Hill Climbing Refinement:

for each node pair (i, j) with $i \neq j$ do

 if there is no edge $(i \rightarrow j)$ in G then

 Add edge $(i \rightarrow j)$ to G .

 if G is acyclic then

 Compute $S_{\text{new}} \leftarrow \text{BayesianScore}(G, D)$.

 if $S_{\text{new}} > S$ then

$S \leftarrow S_{\text{new}}$

 else

 Remove edge $(i \rightarrow j)$

 else

 Remove edge $(i \rightarrow j)$

 else

 Remove edge $(i \rightarrow j)$.

 Compute $S_{\text{new}} \leftarrow \text{BayesianScore}(G, D)$.

 if $S_{\text{new}} > S$ then

$S \leftarrow S_{\text{new}}$

 else

 Re-add edge $(i \rightarrow j)$

Output:

Return the final DAG G with the highest Bayesian score S .

For the K2 algorithm, I did try different modifications, like:

- Taking 10 random node orders, doing K2 on all of them and then picking the graph with the maximum Bayesian Score.
- Picking an edge which reduces the score a small fraction of the time (around 10% of the time).
- Trying a Mutual-Information informed ordering for the nodes, which means that nodes with a higher total Mutual Information will be earlier in the ordering.

Ultimately, I just settled on a regular node ordering for K2 (1, 2, ... n) and did local search afterwards on K2 graph. This seemed to give the best final Bayesian score.

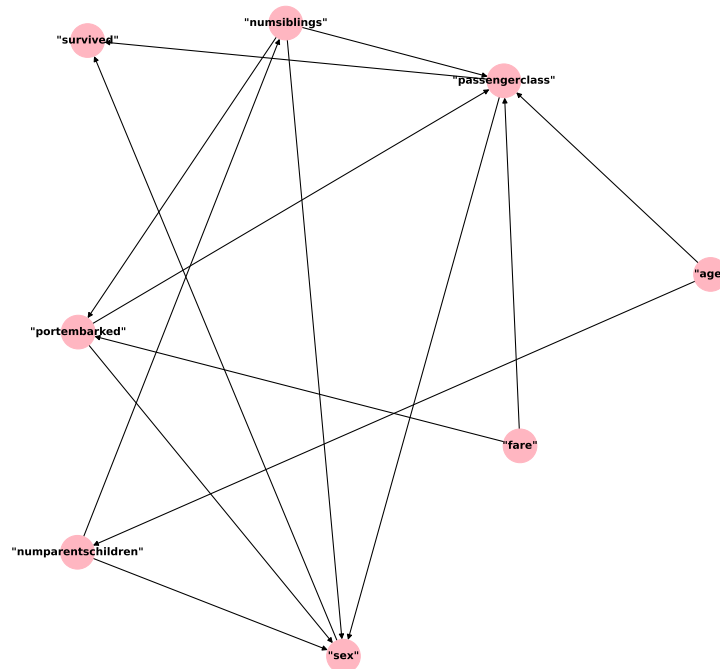
2 Running Times

The running time for various datasets is given below:

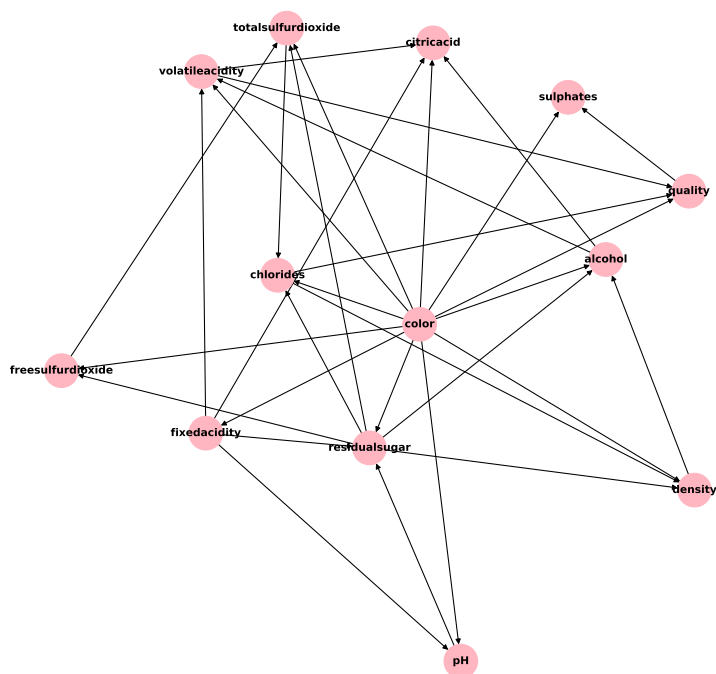
- **small.csv:** 0.25 seconds
- **medium.csv:** 11.07 seconds
- **large.csv:** 1145.86 seconds

3 Visualizing Graphs

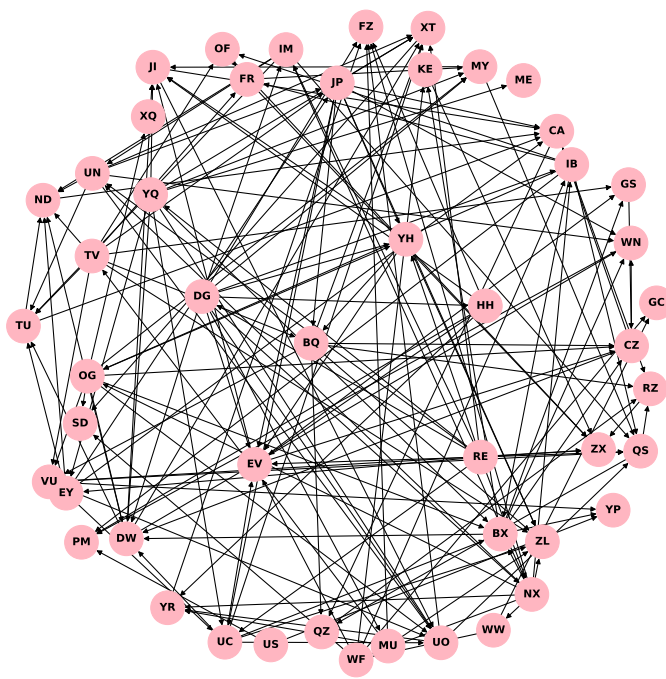
3.1 Final Graph for small.csv



3.2 Final Graph for medium.csv



3.3 Final Graph for large.csv



4 Code

First, we define all our helper functions in `utils.py`:

```
1  # utils.py
2  import numpy as np
3  import math
4  from collections import defaultdict
5  import networkx
6  from sklearn.metrics import mutual_info_score
7
8  def count_configurations(node, parents, data):
9      """
10     Count occurrences of each configuration of a node given its parents.
11
12     Parameters:
13     node (int): The index of the node.
14     parents (list of int): List of parent node indices.
15     data (np.ndarray): 2D numpy array.
16
17     Returns:
18     dict: item
19     """
20     # initialize a dictionary of dictionaries initialized to 0
21     counts = defaultdict(lambda: defaultdict(int))
22     # print("Counting configurations for node:", node, "with parents:", parents)
23     for row in data:
24         parent_values = tuple(row[parents]) if parents else ()
25         node_value = row[node]
26         # print("Row:", row, "Parent values:", parent_values, "Node value:",
27             ↪ node_value)
28         counts[parent_values][node_value] += 1
29     return counts
30
31 def bayesian_score(graph, data):
32     """
33     Calculate the Bayesian score of a given graph structure based on the provided
34     ↪ data.
35
36     Parameters:
37     graph (list of tuples): Each tuple is (node_index, [parent_indices])
38     data (np.ndarray): 2D numpy array of shape (num_samples, num_variables)
39
40     Returns:
41     float: The Bayesian score of the graph.
42     """
43     score = 0.0
44     for node, parents in graph:
45         # print("Node:", node, "Parents:", parents)
46         counts = count_configurations(node, parents, data)
47         # count the number of unique values the node can take
48         unique_node_values = set(data[:, node])
49         # print("Unique values for node {}: {}".format(node, unique_node_values))
50         for parent_values, val_dict in counts.items(): # all key value pairs - here
51             ↪ it is key : (parent values), value : {node val: count}
```

```

50         # print("Parent values:", parent_values, "Value counts:", val_dict)
51         m_ij0 = sum(val_dict.values()) # total count of all values for this
        ↪ parent configuration
52         for m_ijk in val_dict.values():
53             score += math.lgamma(m_ijk + 1) - math.lgamma(1) # Assuming uniform
        ↪ prior : alpha_ijk = 1 for all k
54         score += math.lgamma(len(unique_node_values)) - math.lgamma(m_ij0 +
        ↪ len(unique_node_values))
55
56     return score
57
58 def mutual_info_order(data):
59     """
60     Compute the mutual information between all pairs of variables and return an
        ↪ ordering based on it.
61     We do this to get a better initial ordering for the K2 algorithm.
62
63     Parameters:
64     data (np.ndarray): 2D numpy array of shape (num_samples, num_variables)
65
66     Returns:
67     list: List of variable indices ordered by their mutual information.
68     """
69     num_vars = data.shape[1]
70     mi_matrix = np.zeros((num_vars, num_vars))
71
72     # Calculate mutual information for each pair of variables in given dataset
73     for i in range(num_vars):
74         for j in range(i + 1, num_vars): # mutual information is symmetric
75             mi = mutual_info_score(data[:, i], data[:, j])
76             mi_matrix[i, j] = mi
77             mi_matrix[j, i] = mi
78
79     # Sum mutual information for each variable
80     mi_sums = np.sum(mi_matrix, axis=1)
81     # Get ordering based on ascending mutual information sums
82     order = np.argsort(mi_sums).tolist()
83
84     return order

```

The actual data processing and structure learning algorithm are defined in `project1.py`. We call helper functions from `utils.py` wherever required.

```

1  # project1.py
2  import sys
3  import math
4  import numpy as np
5  from utils import bayesian_score, mutual_info_order
6  import networkx
7  import time
8
9
10 def write_gph(dag, idx2names, filename):
11     with open(filename, 'w') as f:
12         for edge in dag.edges():
13             f.write("{} , {} \n".format(idx2names[edge[0]], idx2names[edge[1]]))

```

```

14
15
16 def compute(infile, outfile):
17     # converting data csv to numpy array
18     data = np.loadtxt(infile, delimiter=',', skiprows=1, dtype=int)
19
20     # mapping variable names to indices to make computation faster (no more dicts)
21     num_vars = data.shape[1]
22     var_names = np.loadtxt(infile, delimiter=',', dtype=str, max_rows=1)
23     name2idx = {var_names[i]: i for i in range(num_vars)}
24     idx2names = {i: var_names[i] for i in range(num_vars)}
25
26     # deciding an order for K2 algorithm
27     order = list(range(num_vars))
28     # initialize graph with no edges
29     dag = networkx.DiGraph()
30     dag.add_nodes_from(range(num_vars))
31     # get the graph for the initial (no edges) structure
32     graph = [(i, list(dag.predecessors(i))) for i in range(num_vars)]
33     current_score = bayesian_score(graph, data)
34
35     # time the algorithm
36     start_time = time.time()
37
38     # K2 algorithm
39     for i in range(num_vars):
40         node = order[i] # for each node in the order, try to add right children
41         for potential_child in order[i+1:]: # only consider nodes that come after it
42             ↪ in the order
43             dag.add_edge(node, potential_child) # add the edge
44             new_graph = [(j, list(dag.predecessors(j))) for j in range(num_vars)]
45             new_score = bayesian_score(new_graph, data)
46             if new_score > current_score: # if score improves, keep the edge and
47                 ↪ update score
48                 current_score = new_score
49             else: # otherwise remove the edge
50                 dag.remove_edge(node, potential_child)
51
52     # now, take this graph and do greedy hill climbing to improve it further
53     for i in order:
54         for j in range(num_vars):
55             if i == j: # self loops not allowed
56                 continue
57             # try adding edge i -> j if it doesn't create a cycle
58             if not dag.has_edge(i, j):
59                 dag.add_edge(i, j)
60                 if networkx.is_directed_acyclic_graph(dag): # only keep if no cycle
61                     new_graph = [(k, list(dag.predecessors(k))) for k in
62                                     ↪ range(num_vars)]
63                     new_score = bayesian_score(new_graph, data)
64                     if new_score > current_score:
65                         current_score = new_score
66                     else:
67                         dag.remove_edge(i, j)
68             else:
69                 dag.remove_edge(i, j)

```

```

67         else:
68             # try removing edge i -> j
69             dag.remove_edge(i, j)
70             new_graph = [(k, list(dag.predecessors(k))) for k in
71                 ↪ range(num_vars)]
72             new_score = bayesian_score(new_graph, data)
73             if new_score > current_score:
74                 current_score = new_score
75             else:
76                 dag.add_edge(i, j)
77
78     end_time = time.time()
79     print("Time taken: {:.2f} seconds".format(end_time - start_time))
80     print("Final score: {}".format(current_score))
81     write_gph(dag, idx2names, outfile)
82     print("Wrote graph to {}".format(outfile))
83
84 def main():
85     if len(sys.argv) != 3:
86         raise Exception("usage: python project1.py <infile>.csv <outfile>.gph")
87
88     inputfilename = sys.argv[1]
89     outputfilename = sys.argv[2]
90     # time the compute function
91     compute(inputfilename, outputfilename)
92
93 if __name__ == '__main__':
94     main()

```
